

12

EUROPEAN PATENT APPLICATION

21 Application number: 88480101.0

51 Int. Cl.⁵: G06F 9/46, G06F 13/16

22 Date of filing: 23.12.88

43 Date of publication of application:
27.06.90 Bulletin 90/26

64 Designated Contracting States:
DE FR GB

71 Applicant: International Business Machines Corporation
Old Orchard Road
Armonk, N.Y. 10504(US)

72 Inventor: Tran-Gia, Phuoc, Prof.Dr.Ing.
Lehrstuhl für Informatik III Universität
Würzburg
Am Hubland, D-8700 Würzburg(DE)
Inventor: Pauporté, André
Lotissement du Colombier
F-06480 La Colle sur Loup(FR)

74 Representative: Möhlen, Wolfgang C.,
Dipl.-Ing.
IBM Corporation Säumerstrasse 4
CH-8803 Rüschlikon(CH)

54 Load balancing technique in shared memory with distributed structure.

57 In a system with several users MU using a common or shared storage comprising a plurality of memory banks MB, a memory load balancing technique is introduced which causes an even distribution of the storage load among the whole set of memory banks.

For each user MU (in its memory interface MI) a table or map is kept reflecting the load status of all memory banks. The assignment of a new storage record is requested only from a memory bank with relatively low load status.

Upon each memory command which is sent from a user to a memory bank (for creating a record and for storing or fetching data), a reply (acknowledgement) is returned to the user. With each such reply, an indication of the current load status of the respective memory bank is returned, and the user's load status table is updated. Thus, no extra data transfers are necessary to keep all users up-to-date.

GLOBAL CONFIGURATION

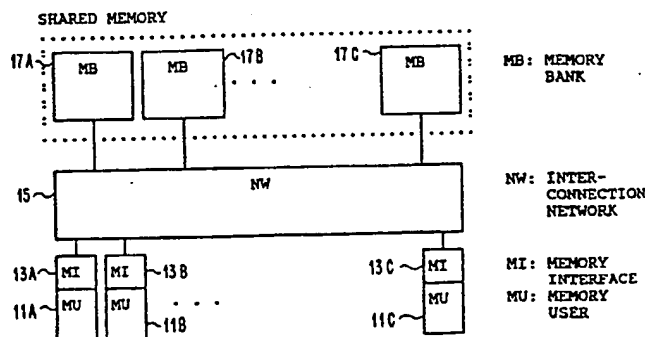


FIG. 1

EP 0 374 337 A1

LOAD BALANCING TECHNIQUE IN SHARED MEMORY WITH DISTRIBUTED STRUCTURE**FIELD OF INVENTION**

Present invention relates to systems where multiple users have access to a common or shared memory consisting of a plurality of memory modules (memory banks) which are separately accessible, and in particular to a method for balancing the load, i.e. the storage utilization among all memory modules.

BACKGROUND

For systems having a plurality of shared resources such as processors, I/O devices, buffers etc. it is of interest to balance the load or utilization between the resources to optimize the performance of the system. Several solutions for balancing the utilization among multiple devices have been suggested, as is reflected in the following list of patents and publications.

U.S. Patent 3,702,006 entitled "Method for Balancing the Utilization of Input/Output Devices". -- U.S. Patent 4,403,286 entitled "Balancing Data-Processing Workloads". -- U.S. Patent 4,633,387 entitled "Load Balancing in a Multiunit System".

Publications in IBM Technical Disclosure Bulletin: -- J.A. McLoughlin: "Load Balancing Buffered Peripheral Subsystem"; Vol. 24, No. 1B (June 1981) pp. 707-709. -- D.C. Cole et al.: "Recalibrating Buffer Peripheral Systems"; Vol. 24, No. 2 (July 1981) pp. 1041-1042. -- G.A. McClain et al.: "Multiple Utilization Threshold Control for I/O"; Vol. 24, No. 3 (Aug. 1981) p. 1406. -- D.W. Bolton et al.: "Combined Real-Virtual Fixed-Size Buffer Mechanism in Shared Storage"; Vol. 27, No. 4B (Sept. 1984) pp. 2655-2660.

In all of these known solutions, either a centrally-kept record is required for holding the information on the load or utilization status, or central processing is provided to obtain the necessary information and to exercise control for load balancing, and furthermore particular transmissions must be performed to collect or distribute the load status and/or control information involved.

The central storage or processing of balancing information can lead to a congestion when many transactions are made simultaneously, and it is detrimental if the respective unit fails, and the requirement for additional data transfers decreases the performance. Furthermore, the modular expansion of a system may not be possible if the centralized updating procedures are designed for a given configuration.

OBJECTS OF THE INVENTION

It is a primary object of the invention to devise a technique for balancing the load among plural storage modules forming a shared storage for several users, in a fully decentralized manner.

It is a further object of this invention to provide a load balancing technique for a shared modular memory, which does not require particular data transfers for the exchange of load status information.

Another object of the invention is to devise a memory load balancing mechanism which allows a modular expansion of the number of memory modules or banks without a necessity for changes or adaptations in the user units, their interfaces and attachments, or in the commonly used units.

SUMMARY OF THE INVENTION

These objects are achieved by a method of shared memory load balancing as defined in Claim 1.

This invention allows balancing of the storage load with a minimum of extra operations and no additional information transfers, by keeping a memory load status map with each user unit, and by updating it in convenient situations. For the updating, the transfer of separate information is avoided by using as carriers the command replies which are transferred frequently from memory banks to user units anyway.

Further features and advantages of the invention will become apparent from the following detailed description of a preferred embodiment in connection with the accompanying drawings.

LIST OF DRAWINGS

Fig.1 is a block diagram of the overall configuration of a multiple user system with shared storage in which the invention is used;

Fig.2 is a schematic representation of the information flow and sequence of steps for a memory access in the system of Fig.1;

5 Fig.3 shows the formats of the commands and command responses which are used for storage operations in the system of Fig.1;

Fig.4 shows the formats of memory control blocks which are stored in the user's local memory and are used for information transfer between the user and its memory interface;

10 Fig.5 is a block diagram of the memory interface (MI) interconnecting a user and the shared memory in the system of Fig.1;

Fig.6 is a flow diagram of the memory command processing in the memory interface logic; and

Fig.7 is a flow diagram of the selection of a memory bank, as done by the memory interface using the load balancing technique of present invention.

15 DETAILED DESCRIPTION

20 1) SYSTEM OVERVIEW

Fig.1 shows the configuration of a system in which present invention finds application. In principle it is a system in which multiple users are connected to multiple memory banks (memory modules) which function as common (or shared) memory for the whole set of users. Such arrangement is used, e.g., in communication controllers.

25 Each memory user MU(11A, 11B, ... 11N) is connected through its own memory interface MI (13A, 13B, ... 13N) to a common interconnection network NW (15). The memory banks MB (17A, 17B, ... 17N) are equal units which are individually connected to the interconnection network.

The interconnection network is a multiple port packet switch which, in response to data packets entered at one of its ports, can route those to the required output port. Such networks are known and are described 30 e.g. in the following publications, and in the literature cited therein. T.Feng: "A Survey of Interconnection Networks"; Computer (IEEE) December 1981, pp.12-30. --V.P.Kumar et al.: "Augmented Shuffle-Exchange Multistage Interconnection Networks"; Computer (IEEE) June 1987, pp.30-40.

The advantage of the whole arrangement is that it can be modularly expanded (as long as the interconnection network is of sufficient capacity and has enough ports). The number of memory banks can 35 be adapted to the number of users and to the storage requirements they have.

The overall flow of commands and control information for memory accesses is schematically shown in the diagram of Fig.2

The memory user MU prepares a memory control word MCW specifying the desired operation and containing parameters. The MCW is transferred to the memory interface MI which processes the command 40 request and sends a command to one of the memory banks MB. The command is executed in the MB (storing of data, reading of data, etc.) and a command response is returned to the MI. The MI prepares status data reflecting the response received from the MB and possibly additional state data reflecting how the memory operation completed, and these status data are transferred in a memory status word MSW to the MU.

45 The term "command" as used in the following description actually designates a request for a memory operation, which is sent from a Memory User/Memory Interface to a Memory Bank (such as CREATE command, PUT command, etc).

The formats of the commands and command responses, of the MCW and of the MSW are detailed later in connection with Figs.3 and 4.

50 The command processing is done by a logic layer (in the MI) which downwards, has an interface to the memory user MU allowing read/write access to the MU local memory, and upwards, has another interface to the interconnection network NW. Details of the MI logic and its operation will be described later in connection with Figs.5 and 6.

One principle assumed for the present system is that the distribution of user data among the storage 55 banks is arbitrary, i.e. there is no fixed association or preassignment between memory banks and users. One problem arising from this principle is that the memory banks may be utilized (loaded) unevenly. Thus, a few memory banks may be completely filled with data while others are not used at all. This may result in congestion, unnecessary waiting times, and loss of parallelism which could be avoided as long as the total

memory space available is sufficient for all the data to be stored. The present invention solves this problem by balancing the load between all the memory banks.

5 2) TECHNIQUE FOR LOAD BALANCING

In present description, the term "load" is used in the following sense: The load of a memory bank is the amount of storage space (expressed in units such as bytes, words, or blocks of 256 bytes) which is actually assigned to records currently existing within the respective memory bank.

The "load status" of a memory bank is the result of a comparison between the current load of the memory bank and a predetermined threshold. This threshold may e.g. be 70% of the total storage space available in the respective memory bank. This load status can be represented by a binary "load status indicator". (If desirable, it would of course be possible to define several levels of load status. Thus, the actual load could be compared to three different thresholds, resulting in four different load status indications which could be represented by two bits.)

The solution of present invention for balancing the load (storage utilization) among all installed memory banks can be summarized as follows:

- a) Keep with each user (e.g. in the memory interface) a table (a map) about the load status of each installed memory bank;
- b) When a user wants to create a new record in the shared memory, select one of the memory banks by interrogating the load status table (to select any low load memory bank if possible);
- c) With each information block flowing back anyway from the memory bank to a user (e.g. acknowledgement of a memory command received from and executed for that user), include some information on the load status of the respective memory bank;
- d) Update the load status table kept with the user (or memory interface) by entering the received load status information.

This novel technique allows an even distribution of the load among all the memory banks with a minimum in extra operations, and requiring no additional data transfers at all because the load status information is carried in information blocks which are transferred anyway.

3) MEMORY INTERFACE FUNCTIONS AND COMMAND FORMATS

The function of the memory interface MI is to provide a standard decentralised memory interface to each memory user MU. It handles the memory commands which are presented by the memory user MU as follows. It selects the appropriate memory bank MB (using the load balancing technique of the present invention), sets up a connection to that MB through the interconnection network NW, builds and sends a command (which may also include data) over that connection, waits for a command response (which in certain cases also includes data from the MB), releases the connection and generates completion status data for the memory user MU.

Each command is composed of a 12-byte header followed by an optional data field. Each command response is composed of an optional data field followed by a 5-byte trailer.

Four different memory commands and associated command responses are provided. Their formats are illustrated in Fig.3, and they are briefly described in the following:

CREATE (format: Cmd only)

This command causes creation of a RECORD (a logical unit of storage) in the shared memory and the assignment of a Logical Record Address (LRA) to it. The LRA is a system unique reference which allows to retrieve a memory record in the shared memory. A memory bank MB will be chosen according to the new load balancing technique presented here.

No memory space is allocated to a new record (as long as no data are to be stored yet).
The command response returns the LRA to the respective memory interface MI.

PUT (format: Cmd, LRA, D, data)

The PUT command causes the data it carries with it to be written into the specified record (LRA) at the specified displacement D. Enough memory space is allocated to the record (when the PUT command is executed) so that the data can be stored.

The command response carries no specific data back to the memory interface MI (except for the Return Code RC contained in each command response).

10 GET (format: Cmd, LRA, D, N)

The Get command causes reading of N data bytes from the specified record (LRA) at the specified displacement D.

The command response carries the N data bytes back to the memory interface MI.

15

RELEASE (format: Cmd, LRA)

The Release command causes cancelling of the specified record (LRA) and releasing of the associated memory space.

The command response contains no specific data (except for the Return Code).

Return Code Structure and Load Status Indicator:

25

The Return Code (RC) which is part of each command response returned by a memory bank is an 8-bit byte which has the following structure: The first four bits RRRR are a specific code which indicates the completion status of the respective memory command. The fifth bit L is a load status indicator (LSI) whose binary value indicates the load status of the respective memory bank: A 1 stands for low load (below a given threshold) and a 0 indicates high load.

30 The last three bits of the RC are unused in present example. However, they could be utilized to distinguish between several levels of storage load: With a two-bit load status indicator L (instead of the one bit mentioned above), one can distinguish four levels of utilization or load.

35 The LSI returned in the command response is generated by the MB based on the comparison of two quantities: The "free memory space" and the "free memory space threshold". The free memory space is the amount of memory locations in the MB which are not assigned to any record. That quantity is maintained by the MB control during command execution. The Create and Get commands which do not consume free memory space do not alter that quantity. The Put commands which consume memory space decrease the free memory space. The Release command will increase the free memory space. The free memory space threshold is a value which has been initialized in each MB.

40 The operation of the MBs will not be further described as it is not the subject of present invention.

4) MEMORY USER AND MEMORY CONTROL BLOCKS

45

The memory user MU can be viewed as a microprocessor with its main store which is called the MU local memory. The memory commands move data between the shared memory (the total set of all memory banks MB) and the MU local memory.

50 Information exchange between a memory user MU and its memory interface MI is effected through memory control blocks the formats of which are shown in Fig.4.

A memory command is presented by the memory user MU as a control block that is called "Memory Command Word" (MCW) and is set up in the local memory of the memory user MU. The MCW contains control information specifying the type of the memory command (create, get, put, release), the Logical Record Address (LRA), the data displacement (D) within the record, the data count (N), and the address of the data area within the MU local memory (local memory pointer LMP).

55 The completion status data is presented by the memory interface MI to the memory user MU as a control block in the MU local memory: the "Memory Status Word" (MSW). The status data contain information (CC) about whether the command execution was satisfactory or not. In the case of a Create

command, the MSW also contains the Logical Record Address (LRA) of the created record.

5) MEMORY INTERFACE LOGIC

A block diagram of the memory interface logic is shown in Fig.5. Its elements and structure are as follows:

The MI basically comprises a central portion 31, an interconnect network (NW) port 33, a DMA portion or user port 35, and a control section 37.

The central portion 31 includes an ALU 41, an accumulator register (ACCU) 43, local storage LS 45 and local storage address register LSAR 47. An A bus 49 and a B bus 51 are connected to the inputs of the ALU, and a result bus 53 is connected to the output of the ALU. Input data to the LS are transferred on lines 55, 57, and 59, and output data from the LS are transferred on line 61.

NW port 33 comprises four registers and a counter. NW Data In (NWDI) register 63 is connected between LS output line 61 and a two-way NW data bus 65. NW Data Out (NWDO) register 67 is connected between LS data input line 55 and the NW data bus 65. NW Control In (NWCi) register 69 is connected between ALU result bus 53 and an NW control input line 71. NW Control Out (NWCO) register 73 is connected between A bus 49 and an NW output control line 75. An NW Counter (NWCT) 77 receives its input from ALU result bus 53. Its contents can be incremented (or decremented) by one unit when a byte was received from (or furnished to) the network, and it can be tested for zero.

User port 35 includes three registers and a counter. MU Data IN (MUDI) register 79 is connected between LS input line 57 and a two-way user data bus 81. MU Data Out (MUDO) register 83 is connected between LS output line 61 and the user data bus 81. MU Address (MUA) register 85 is connected between ALU output bus 53 and a user address line 87. An MU Counter (MUCT) 89 receives its input from ALU result bus 53. The contents can be decreased by one unit when a byte was received from (or furnished to) the user, and it can be tested for zero.

Control section 37 furnishes the following data or control signals:

IF =	Immediate Field (LS Operand or LS Address)
ALUOP =	ALU Operation code (Add, Subtract, And, Or, Xor, Rotate,...)
ABS =	A Bus Source
BBS =	B Bus Source
LSOP =	LS Operation code (Read, Write, NOOP)
LSA =	LS Address Control
MISC =	Miscellaneous controls to write registers, start DMA, etc.

The MI logic has the following capabilities:

1- The MI Logic can perform arithmetic and boolean operations on two operands, one being an immediate value (IF) generated by the control section, and the second one being the contents of an LS location or of a control register.

2- The MI logic has read/write access to the MU local memory, the local memory address and the byte count being respectively specified in the MU Address Register (MUA) and the MU Count Register (MUCT). The source of the data during write (or the destination during read) is the LS addressed by the LS Address Register (LSAR).

3- The MI logic can establish a connection to a memory bank MB through the interconnection network NW by loading a connection request and the MB address (identifier) into the "NW Control In Register" (NWCi) and waiting for a connection grant to be returned in the "NW Control Out Register" (NWCO). The connection can be released in a similar way by loading a connection release request into the NWCi and waiting for a grant to be returned in the NWCO. The interconnection network NW will not be further described as it is not the subject of present invention. (Publications describing interconnection networks have been mentioned in section 1 above.) Once the connection is established to a memory bank MB, the MI logic can make a two way data transfer over that connection: assuming the LSAR is initialised to value A and the NW Count Register (NWCT) to value N, the N bytes located at LS(A) are read out and sent to the MB. This inbound data transfer terminates with NWCT=0 and LSAR=A+N. The MB is then supposed to return a string of bytes which are received and stored into the LS starting at address A+N. The number of received bytes is available in the NWCT at the completion of the outbound transfer. The

commands and command responses are exchanged between the memory interface MI and a memory bank MB in that way.

5 6) COMMAND PROCESSING IN THE MEMORY INTERFACE LOGIC

The flow diagram of Fig.6 illustrates the steps which are executed during command processing in the memory interface (MI) logic. In particular, it shows how the memory utilization (memory load status) information which is an important part of the load balancing technique of present invention, is handled during the processing of memory commands.

Step A: The MI reads the memory control word MCW from the user's local memory (using DMA techniques) and stores the MCW in its own local storage LS. It then decodes the command. Depending on the fact whether the command is a CREATE command or not, one of two different steps is taken next.

Step B1 (Cmd = CREATE): One of the memory banks is selected according to the load balancing technique (selection described in Section 7 in connection with Fig.7). The CMB variable designating the current memory bank (explained below) is set up.

Step B2 (Cmd = other than CREATE): The memory bank MB to be used for this command is determined by the LRA field of the MCW. Cf. remarks in Section 7b below. The CMB variable (designating current memory bank) is set up.

Step C: The MI establishes a connection through the interconnection network NW to the memory bank designated by the current CMB. A command is set up in local storage LS (no details given here because this is not essential to the load balancing technique). The starting address of the command in LS is entered into LSAR, and the byte count of the command is loaded into register NWCT. Then the command is sent to the selected memory bank MB. Returning of a command response from the MB is awaited.

Step D: The command response, when received in the MI, is processed (no details given here because this is not essential to the load balancing technique). The L bit of the Return Code is retained as load status indicator LSI.

Step E: The LSI bit is tested to determine the load status of the MB just used. Depending on the result, one of two different steps is taken next.

Step F1 (LSI = 1): The selected MB is in low load state. The LL variable (identifying all low load memory banks, explained below) and the CMB variable (designating the current memory bank) are combined by an OR operation to update LL so that it includes the last-used memory bank (which just reported its load status through the command response) to be in low load state.

Step F2 (LSI = 0): The selected MB is in high load state. The LL variable (identifying all low load memory banks) and the CMB variable (designating the current memory bank) are combined in an AND operation to update LL so that it reflects the high load status for the respective MB.

Step G: Now the MI sets up the memory status word MSW in its local storage LS (using the information received with the command response). It then loads the LS address where the MSW starts into the LSAR and the byte count of the MSW into the counter MUCT, and copies the MSW from the LS into the user's local memory. Finally, the connection through the interconnection network to the MB is released.

Boolean Variables in MI for Memory Bank Selection and Load Status Mapping:

The process uses four boolean variables CMB, RR, TL, and LL which reside in the local store (LS) of the memory interface MI and which are defined below. These four variables designate the memory bank MB currently in use by the respective MI, the Round Robin variable, the total list of all installed memory banks, and the list of memory banks which are presently in a low load status (i.e. which can take more CREATE requests). It is assumed for the examples given below that the maximum number of memory banks is eight. If more than eight memory banks would be provided, these three variables must then comprise two (or more) bytes instead of only one byte as assumed below.

Current Memory Bank (CMB):

8 bits variable numbered 0 to 7.

All bits are zeros except one whose position indicates the selected memory bank MB involved in the

current command.

ex: CMB = 00100000

In this example, MB number 2 is being selected.

5

Round Robin Variable (RR):

8 bits variable numbered 0 to 7.

10 All bits are zeros except one whose position indicates the memory bank MB selected during the last Create command.

ex: RR = 00010000

In this example, MB number 3 was selected during the last Create command issued by the MI.

15

Total List (TL):

8 bits variable numbered 0 to 7.

20 Bit positions in the TL are mapped to the installed memory banks MB.

ex: TL = 11110000

In this example, four MBs are installed in positions 0, 1, 2 and 3.

25 Low Load List (LL)

8 bits variable numbered 0 to 7.

Meaningful bits correspond to installed memory banks MB as specified in TL.

30 A meaningful bit at one indicates that the associated MB is in low load state.

A meaningful bit at zero indicates that the associated MB is in high load state.

ex: LL = xxxx0000

In this example, x = 0 if the associated MB is in high load state and x = 1 if the associated MB is in low load state.

35

7) SELECTION OF A MEMORY BANK BY THE MEMORY INTERFACE

40 There are two different situations when a memory interface MI has to select a memory bank MB for executing a memory command: (a) For a CREATE command, a new memory bank has to be selected according to the load balancing technique of present invention. (b) For all other memory commands the memory bank is known (from the LRA parameter) and the MI can just identify this MB.

45 a) Selection of an MB according to the load balancing technique:

50 As was already explained in the beginning, the selection of a memory bank MB for creating a new record is done by the load balancing technique which utilizes the knowledge about memory load status which is available in each memory interface MI: If there is any low-load MB, the MI selects the next one in a cyclic sequence which starts with the MB that was selected during the last Create command issued by the respective MI (there is an indicator RR in the MI which identifies this last-used MB). If there is no more low-load MB (i.e. all MBs are already heavily loaded), the MI must make a selection among all installed MBs and does this also by the same round-robin mechanism.

55 The operations to be executed by the MI for selection of an MB are shown in the flow diagram of Fig.7. First of all, a test is made whether the set LL (low-load MBs) is empty, by testing whether LL is all zeros. (1) If LL is not all zeros, i.e. there is at least one MB with low-load status, the RR value is increased by 1 (modulo m, the max number of MBs). This is done by a ROTATE1 ALU operation. Then the variables RR and LL are ANDed. If the result is all zeros, the new RR did not hit a low-load MB and the step must be

repeated. If the result is not all zeros, the new RR hit a low load MB and thus can be used for selecting the MB to execute the CREATE command. (2) If LL is all zeros, there is no low-loaded MB and any MB can be used. Also here, the current RR is increased by 1 (modulo m). Then, the new RR and the variable TL are ANDed. If the result is all zeros the new RR did not hit an available MB and the step must be repeated. If
 5 the result is not all zeros the new RR hit an available MB and can be used for selecting the next MB to be accessed for creating a new record.

The selected MB is identified by the contents of RR. This RR contents is now loaded into CMB for further processing.

10

b) Selection of an MB when the LRA is known:

The LRA which is returned by the MB on a CREATE command is a 32 bits identifier in which bits 8, 9 and 10 contain a code specifying the MB (000 identifies MB number 0; 001 identifies MB number 1; etc.).

15 Therefore, the MI logic knows which MB it has to establish a connection to.

Selecting an MB for commands other than CREATE consists in identifying the LRA field in the MCW (which has just been copied in LS), extracting bits 8 to 10 out of it, translating this three bits code into the eight bits format used in the CMB and loading the resulting value into the CMB.

20

Claims

1. Method for evenly distributing the storage load among a plurality of storage units which are commonly used by a plurality of user units, in a system in which a storage operation is caused by a
 25 command sent from a user unit to a storage unit, and in which a command response is returned from a storage unit to a user unit in response to each command, said method comprising the steps of:

- inserting into each command response returned from a storage unit to a user unit, a load status indicator reflecting the storage utilization status of the respective storage unit;
- establishing, with each user unit, a load status table reflecting the storage load status of all installed
 30 storage units of the system;
- updating said load status table in response to each load status indicator received with a command response; and
- prior to sending a command from a user unit for requesting the creation of a new storage record in any storage unit, interrogating said load status table and selecting a storage for which a low load status is
 35 indicated in said table, as recipient for said command requesting a storage record creation.

2. A method in accordance with Claim 1, in which

- in said inserting step, a binary character which indicates either a high load status or a low load status is inserted into the command response; and
- in said establishing step, a load status table is established which comprises for each storage unit installed
 40 in the system at least one binary character position for indicating the load status of the respective storage unit.

3. A method in accordance with Claim 2, comprising the following additional steps:

- defining for each storage unit installed in the system, as load threshold a particular percentage of the total storage capacity of that storage unit; and
- 45 - for determining the load status indicator to be inserted into a command response, comparing the amount of storage space currently assigned to storage records in the respective storage unit, to said load threshold.

4. A method in accordance with Claim 1 or 2, comprising the following additional steps:

- defining a cyclic sequence of all installed storage units;
- storing with each user unit an identification of the storage unit which was selected as recipient of the most
 50 recently issued command for requesting a storage record creation; and
- in said load status table interrogating and storage unit selecting step, selecting from plural storage units having a low load status that storage unit which is the next in said cyclic sequence after the one whose identification was stored.

5. A method in accordance with Claim 1 or 2, comprising the following additional steps:

- 55 - storing in each user unit after sending to a selected storage unit a command requesting a storage record creation, the identification of said storage unit as last-selected storage unit; and
- in said load status table interrogating and storage unit selecting step, if the load status table does not indicate any storage unit having a low load status, selecting that storage unit which is the next in said cyclic

sequence after the last-selected storage unit.

6. In a system comprising several user units, a set of plural storage units commonly usable by said user units, and an interconnection network interconnecting said user units and said storage units; in which system user units send commands to storage units for storage access operations and receive a command response for each command;

5 a method of evenly distributing the storage load throughout all installed storage units, characterized by the following steps:

- establishing with each user unit a storage load status map reflecting the current storage space utilization of each installed storage unit;
- 10 - establishing with each storage unit a load status indicator reflecting its current status of storage space utilization;
- in a user unit, for requesting creation of a storage record, selecting one of the installed storage units in dependence of the contents of the storage load status map and sending a storage record Create command to the selected storage unit;
- 15 - in the selected storage unit, upon reception of the Create command, creating a storage record and allocating a record identifier to it; and returning a command response to the requesting user unit, containing the identifier of the created storage record, and said load status indicator; and
- in the requesting user unit, upon reception of the Create command response, updating the storage load status map in response to the received load status indicator.

20 7. A method in accordance with Claim 6, characterized by the following additional steps:

- in a user unit, for requesting data to be stored into a specified record, selecting the storage unit which contains the specified record and sending a Put command containing said data to the selected storage unit;
- in the selected storage unit, upon reception of a Put command, storing the data into the specified record, allocating memory space for those data if necessary, updating the load status indicator accordingly, and
- 25 returning a command response to the requesting user unit, containing said load status indicator; and
- in the requesting user unit, upon reception of a Put command response, updating the storage load status map in response to the received load status indicator.

8. A method in accordance with Claim 6, characterized by the following additional steps:

- in a user unit, for requesting release of a specified record, selecting the storage unit which contains the
- 30 specified record and sending a Release command to the selected storage unit;
- in the selected storage unit, upon reception of a Release command, cancelling the specified record, deallocating any memory space used for the specified record, updating the load status indicator accordingly, and returning a command response to the requesting user unit, containing said load status indicator; and
- 35 - in the requesting user unit, upon reception of a Release command response, updating the storage load status map in response to the received load status indicator.

40

45

50

55

GLOBAL CONFIGURATION

FIG. 1

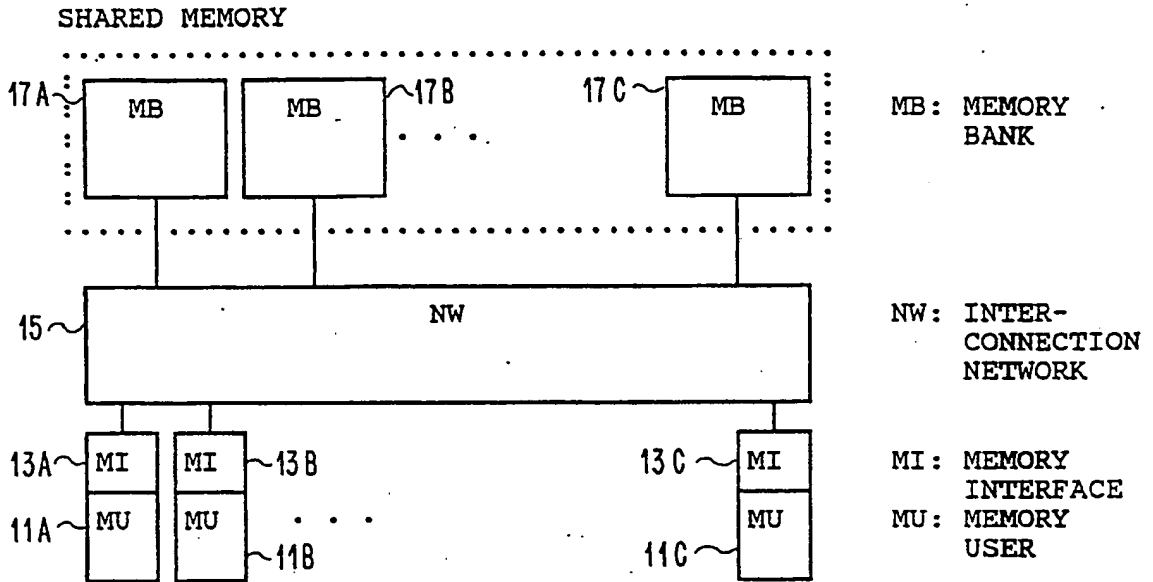
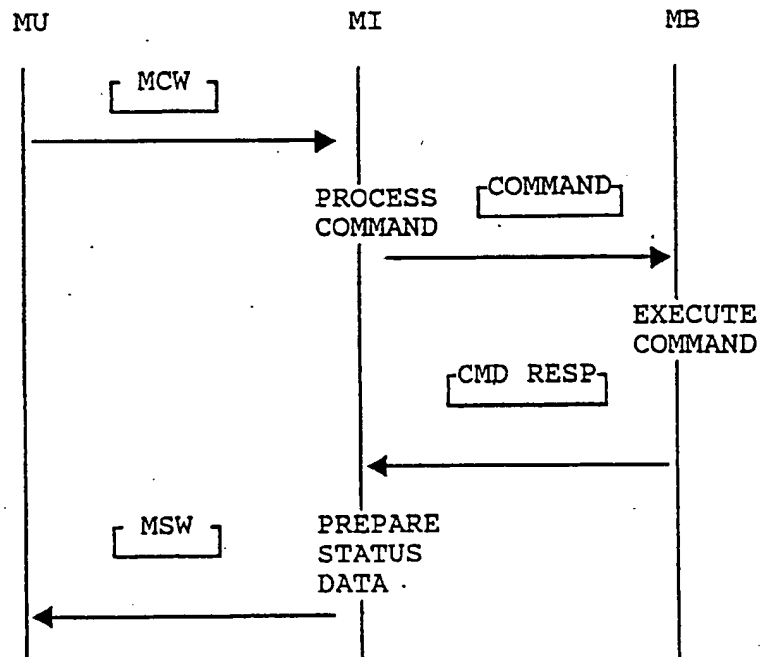
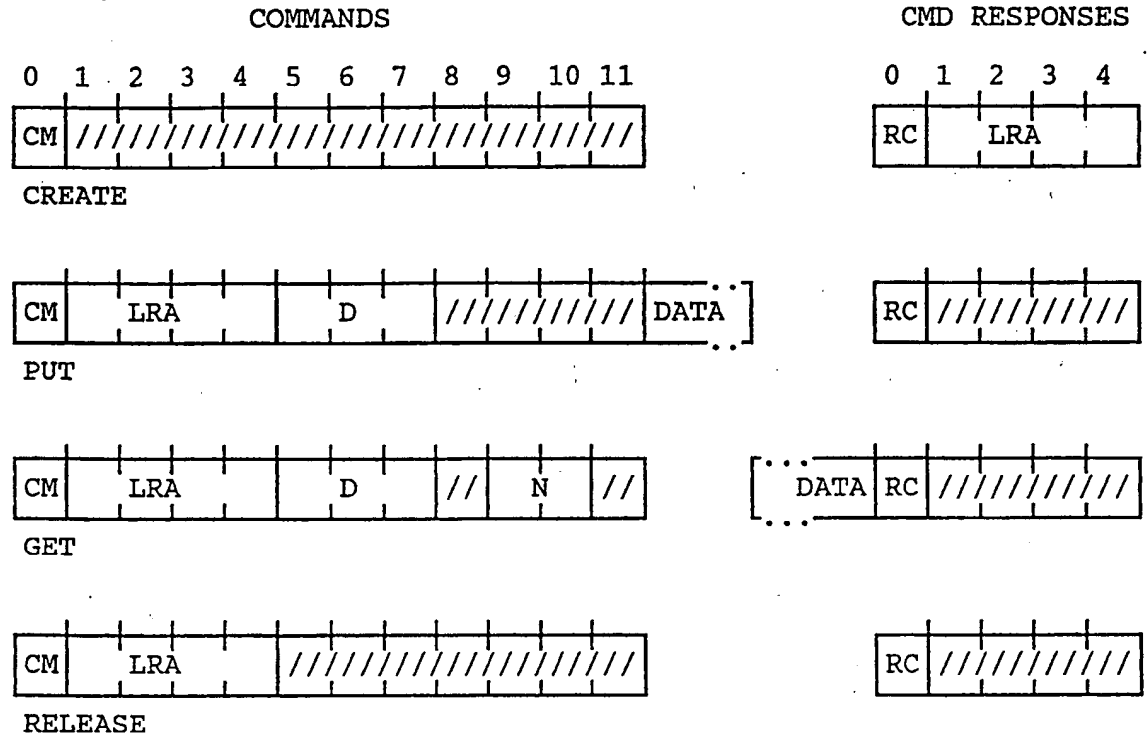


FIG. 2



OVERALL
COMMAND
PROCESSING

FIG. 3 COMMAND FORMATS AND COMMAND RESPONSE FORMATS



CM : Memory Command Code
 LRA : Logical Record Address
 D : Displacement
 N : Byte Count
 /// : Unused

RC : Return Code

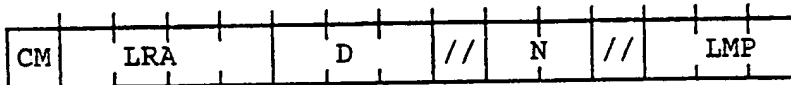
Return Code format: RRRRLxxx (8 bits)

Where RRRR = Command Return Code
 L = Load Status Indicator (LSI)
 xxx = Unused

FIG. 4

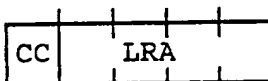
MCW AND MSW FORMATS

MCW



CM: Command Code.
 LRA: Logical Record Address. Unused for the create command.
 D: Displacement. Unused for create and release commands.
 N: Data Byte Count. Used for the get command only.
 LMP: Local Memory Pointer. Used for the get and put commands.
 Points to the come-from or go-to data area in the Memory User Local Memory.

MSW



CC: Command Completion Code.
 LRA: Used in the create command to return the Logical Record Address to the Memory User.

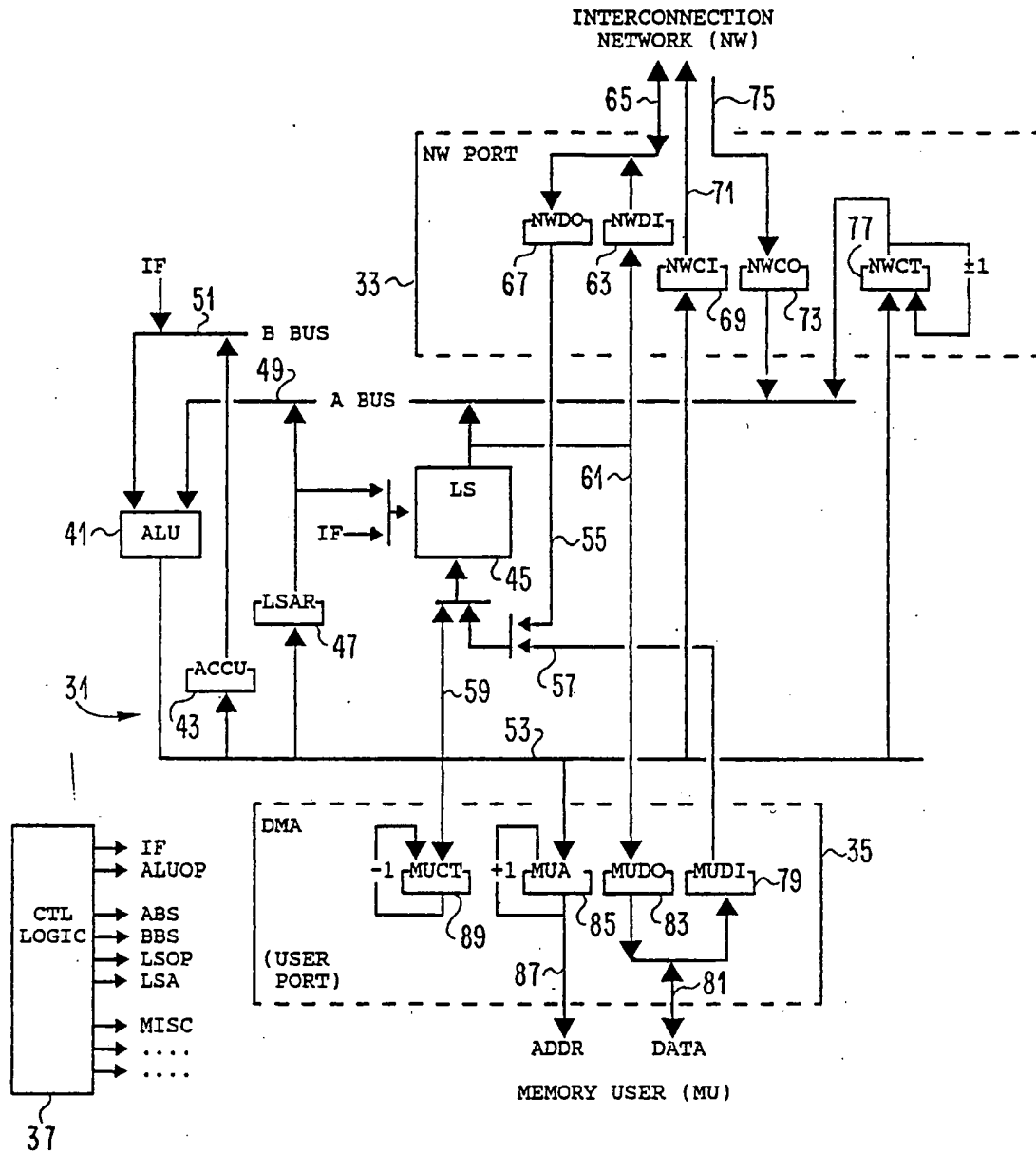
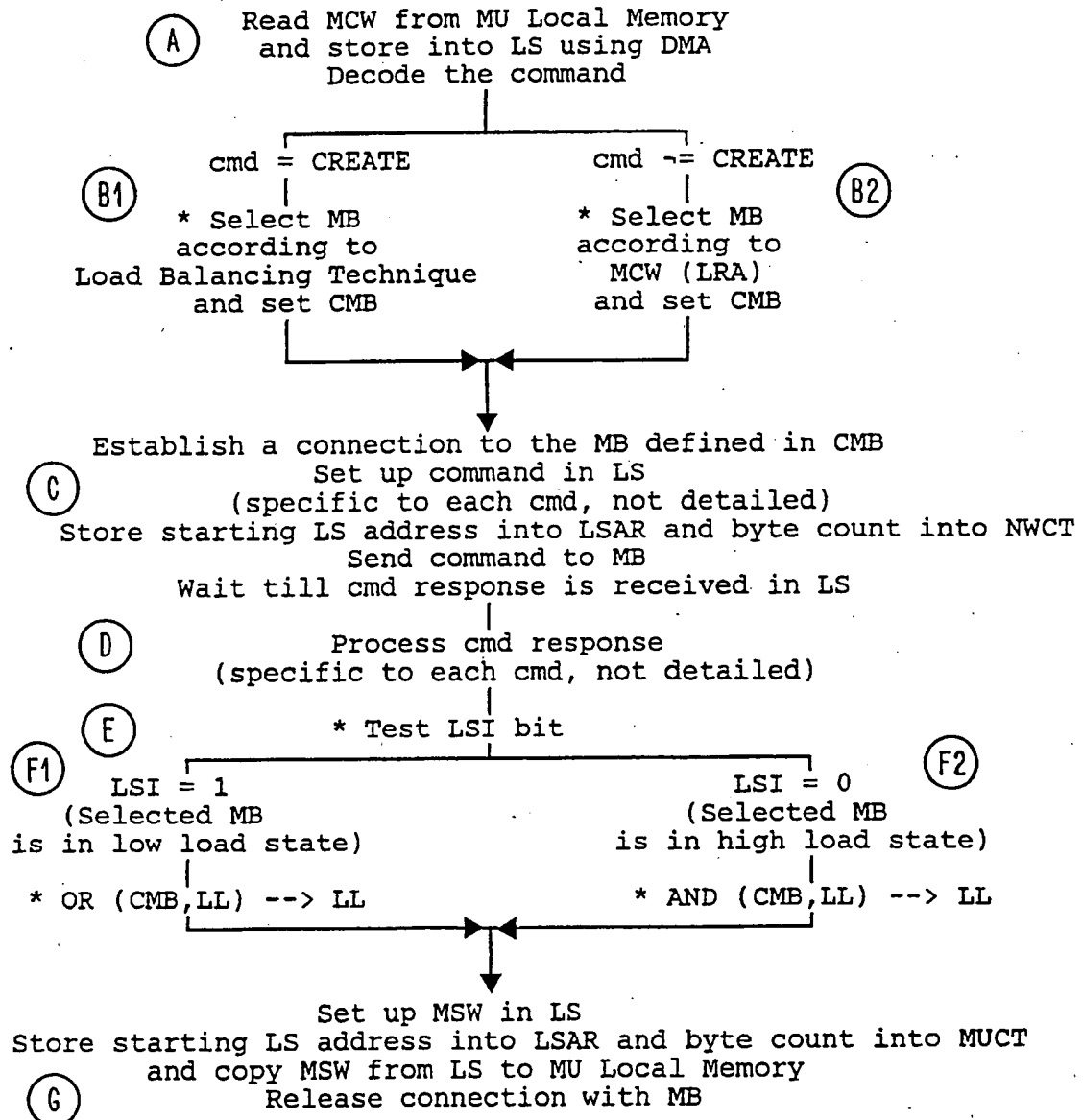


FIG. 5 MEMORY INTERFACE LOGIC

FIG. 6

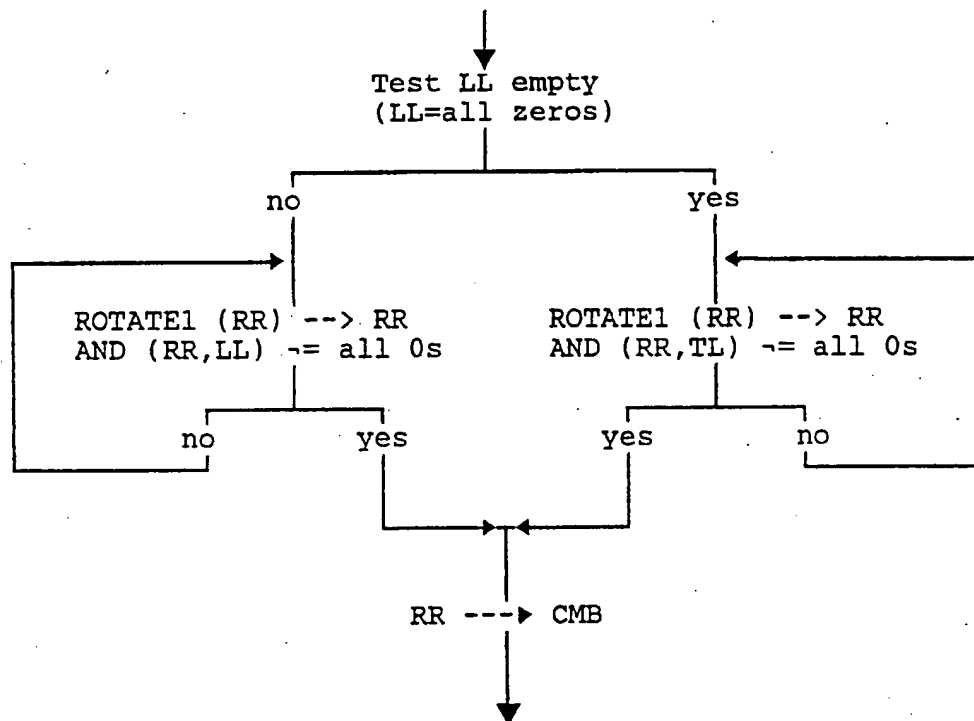
COMMAND PROCESSING IN MI LOGIC



* identifies functions of MI directly associated with
the load balancing technique

FIG. 7

SELECTION OF AN MB WITH LOAD BALANCING



The ROTATE1 ALU operation rotates the operand one bit position to the right.



DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int. Cl.5)
A	PATENT ABSTRACTS OF JAPAN, vol. 8, no. 10 (P-248)[1447], 18th January 1984; & JP-A-58 169 649 (FUJITSU K.K.) 06-10-1983 * Abstract * ---	1,2	G 06 F 9/46 G 06 F 13/16
A	US-A-4 393 459 (HUNTLEY) * Column 2, lines 3-8; column 3, lines 39-63; column 12, lines 43-50; figures 1,8 * ---	1,6-8	
A	EP-A-0 234 803 (TERADATA) * Page 5, lines 24-28; page 6, lines 1,2,23-28; page 7, lines 5-9; page 17, lines 3-20; figure 1 * ---	1,6-8	
A	IBM TECHNICAL DISCLOSURE BULLETIN, vol. 29, no. 3, August 1986, pages 1120,1121, New York, US: "Data set load-balancing algorithm" * Page 1120, lines 11-18; page 1121, lines 8-11 * ---	4,5	TECHNICAL FIELDS SEARCHED (Int. Cl.5)
A	IBM TECHNICAL DISCLOSURE BULLETIN, vol. 19, no. 6, November 1976, pages 2159-2167, New York, US; L. WYMAN: "Intermemory communications facility" * Page 2161, lines 19-41; page 2162, lines 9-25; page 2163, lines 23-41; page 2164, lines 9-20; page 2165, lines 11-30; page 2167 * -----	1,6-8	G 06 F 9/46 G 06 F 13/16 G 06 F 15/16
The present search report has been drawn up for all claims			
Place of search THE HAGUE		Date of completion of the search 29-08-1989	Examiner DHEERE R.F.B.M.
CATEGORY OF CITED DOCUMENTS X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons ----- & : member of the same patent family, corresponding document			